

Bare metal K8s, или Туда и обратно

История Quadcode

Устинов Илья
Kubernetes Platform



HighLoad⁺⁺
2022



quadcode

Про компанию



Quadcode — международная продуктовая компания, которая создаёт программное обеспечение для торговой и инвестиционной индустрии.

≈50M

Пользователей

17

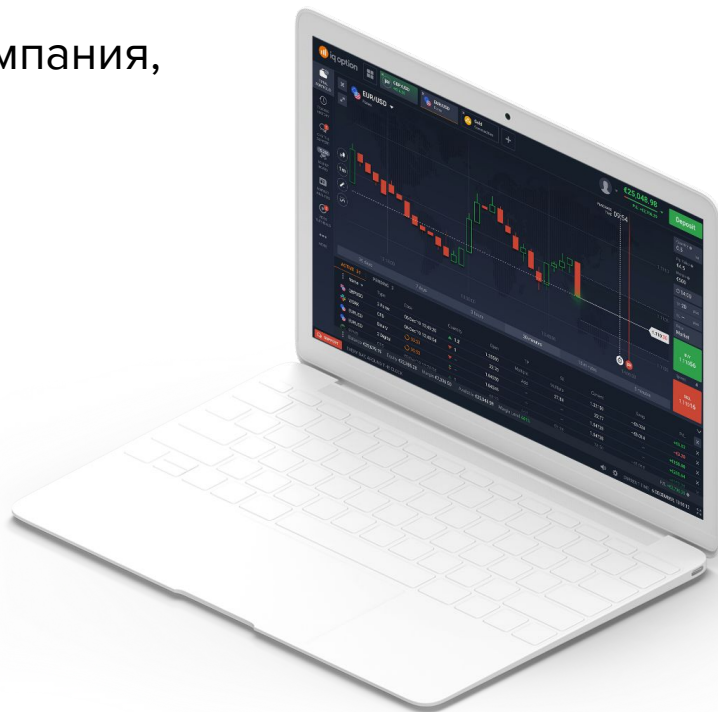
Языков

≈20M

Установок приложений

218+

Стран



Моя команда



Юра
Team Lead



Игорь
DevOps



Кирилл
DevOps



Рома
DevOps



Антон
Product
Manager



Илья
DevOps



Костя
Product
Owner

И принимавшие участие ранее



Вова
DevOps



Саша
Go
Developer



Денис
DevOps



Сергей
DevOps



Даша
Dev



Вова
DevOps



Паша
Dev



Миша
DevOps



Стас
DevOps

Что будет в докладе

Ретроспектива жизни Bare Metal-кластера, а также*:

- Как докатиться до 4 часов при скейле ноды.
- Как сократить это время до 3 минут, когда уже докатились.
- Как AWS, GCE, DO собирают свои кластеры, и что можно перенять у них в Bare Metal.



* возможны флешбэки

Про наши кластеры

3 k8s-кластера, по одному в окружениях:

Infra, Preprod, Prod

	Prod	Preprod	Infra
Количество нод	97 nodes	53 nodes	22 nodes
Каждая нода	40 CPU/256 GB	40 CPU/256 GB	40 CPU/256 GB
Общая ёмкость	3880 CPU, 24.25 TB	2120 CPU, 13.25 TB	880 CPU, 5.5 TB

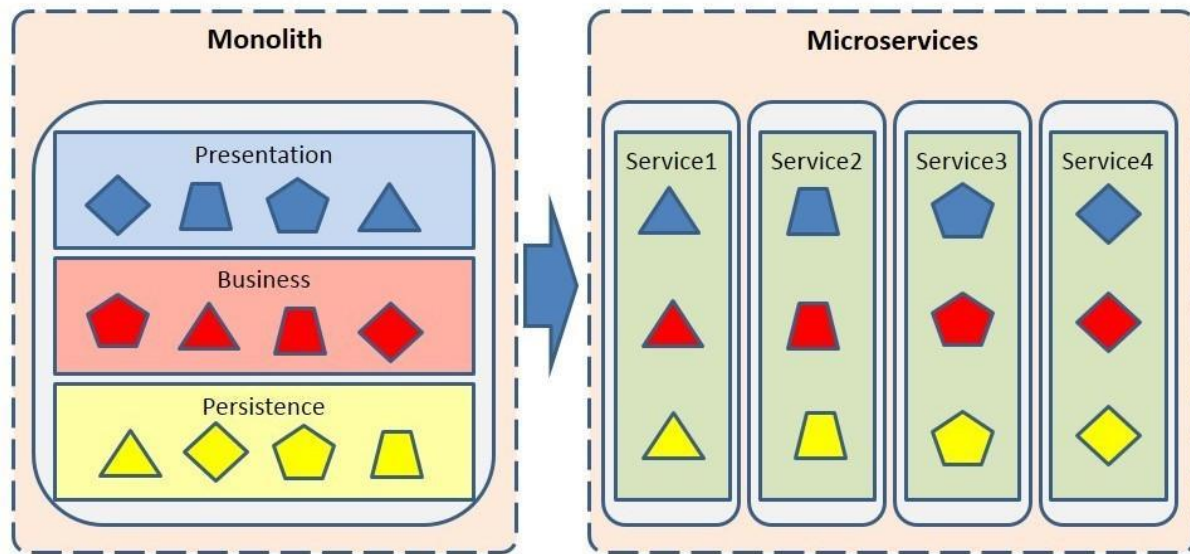
Предпосылки. До 2017 года

Предпосылки



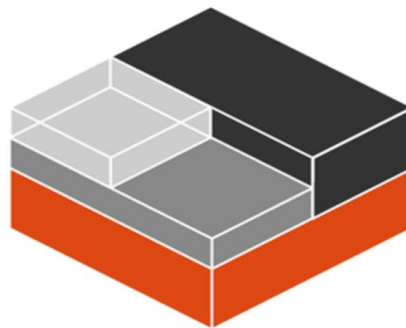
Предпосылки. Архитектура

Монолит → Микросервисы



Предпосылки. Архитектура

Bare Metal + LXC Containers

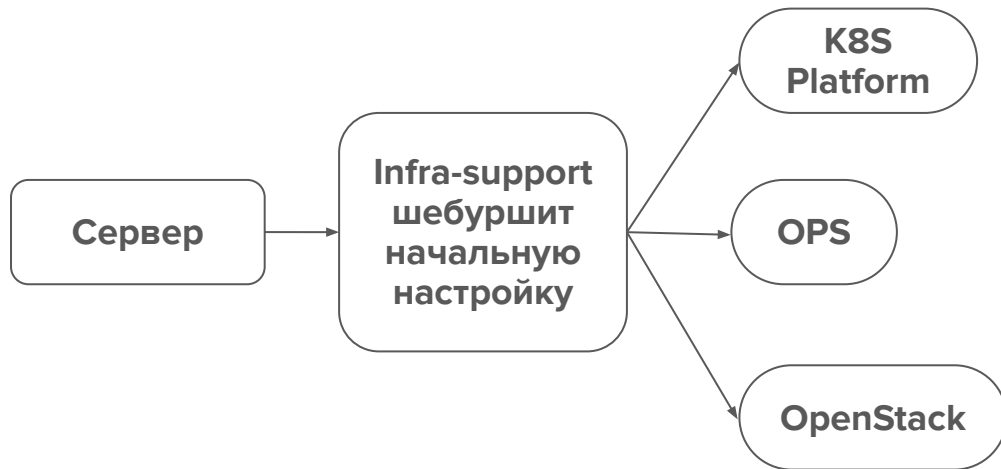


Предпосылки. Архитектура

Отдельная команда эксплуатации. Процесс заказа серверов.

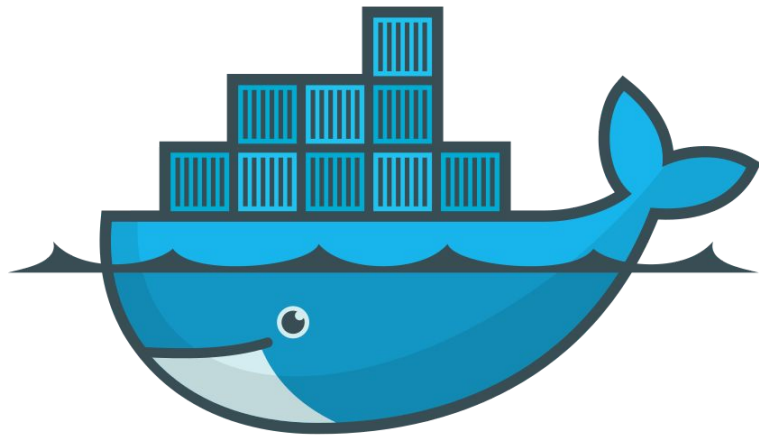


мы настроили



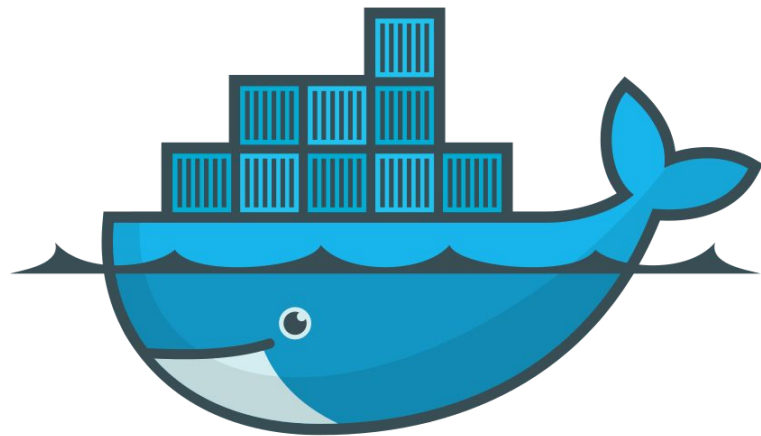
Предпосылки. Архитектура

Docker для dev- и test-окружений

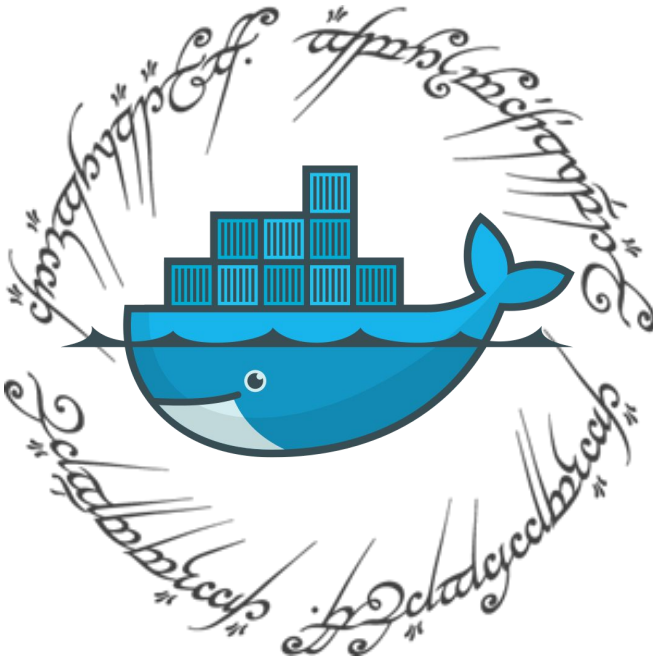


Предпосылки. Итого

- Модное на тот момент желание переехать на микросервисы.
- Много разного неудобного рантайма: вроде LXC, вроде где-то просто настроенные сервера.

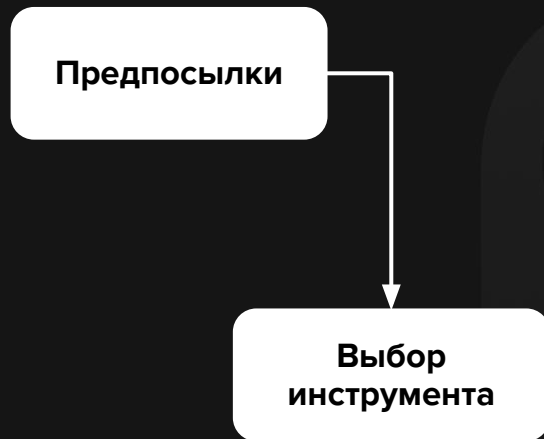


Хотим один рантайм, чтобы править всеми!



Выбор инструмента оркестрации.

Q2 2017



Выбор инструмента.

Self-hosted или облако

Self-hosted:

- Можно потрогать.
- Требования безопасности и бизнеса:
 - Всё должно быть во внутреннем сетевом периметре.
 - Низкая latency.

Облако:

- Нельзя потрогать.
- Ходят слухи, что управлять облаком проще.
- Требования безопасности и бизнеса:
 - Всё должно быть во внутреннем сетевом периметре — нужно идти к пос'ам и узнавать, сколько стоит.
 - Низкая latency — нужно идти к пос'ам.

Выбор инструмента.

Ожидаемая стоимость кластера, просто

Железо:

- Мастера и рефлексоры поменьше 8 CPU/64 GB — €190 в месяц × 6.
- Ноды побольше 40 CPU/256 GB — стандарт поставки дата-центра €585 в месяц × 5 для начала.

€55 800 в год

Облако:

- Amazon EKS — €75 в месяц.
- Amazon EC2 — €6603 в месяц.

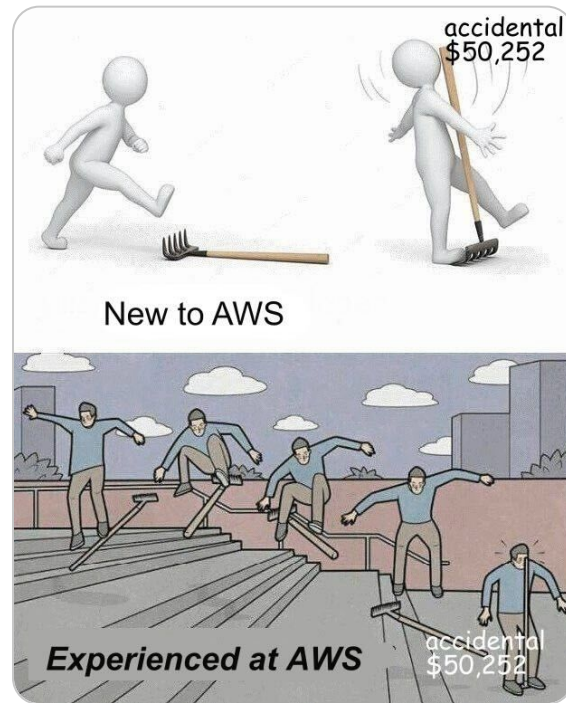
€78 132 в год

Выбор инструмента.

Ожидаемая стоимость кластера

Дополнительная инфраструктура облака —
€78 132 плюс стоимость:

- VPC.
- TGW.
- Direct connect.
- ALB, ELB, NLB, NLP.
- VPC peering.
- ETC.



Выбор инструмента.

Ожидаемая стоимость кластера, сложно.

Экспертиза, люди, процессы



Выбор инструмента. Решения для Bare Metal

Если bare metal, то какой?



- Поднимает кластер в облаке на EC2



Kubernetes The Hard Way
by Kelsey Hightower

- Очень подробный
- Идеально, чтобы учиться
- Нужно автоматизировать

Выбор инструмента.

Kubespary:

- Написан на Ansible.
- Позволяет настраивать компоненты.
- Подготавливает окружение.
- Разрабатывался Mirantis.

Внедрение.

Q3 2017



Внедрение.

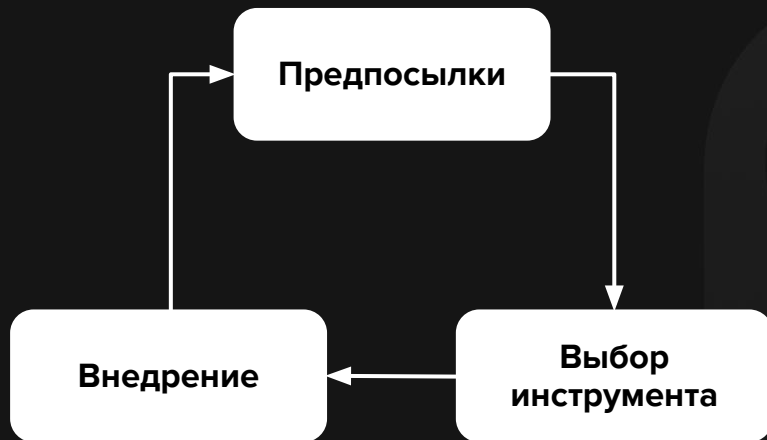
Что включила в себя эксплуатация

Стандартные операции:

- Раскатка кластера.
- Скейл кластера.
- Разработка процесса для ввода приложений в k8s, авторизации, рекомендуемых пайплайнов деплоя, etc.
- Апгрейд софта ноды: Docker, Kernel, etc.
- Апгрейд версии кластера.
- Обновления сертификатов k8s, CNI, etc.

Размер кластера:
5 → 20 нод

Предпосылки рефакторинга. 2017-2019



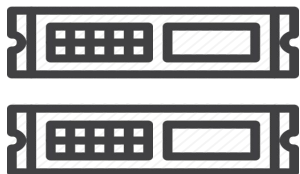
Перерыв на cool story



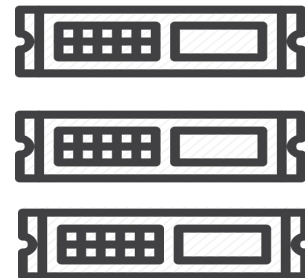
Предпосылки. Проблемы

Физический переезд Master-серверов

Старый DC



Новый DC



Предпосылки. Рост времени выполнения стандартных операций

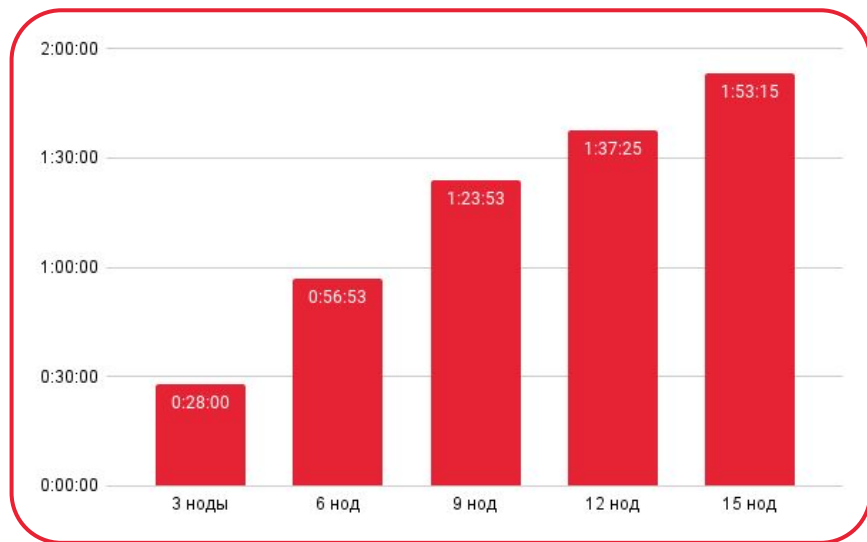
Стандартные операции:

- Раскатка кластера: 0,5 часа → 4 часа.
- Скейл кластера: 0,5 часа → 4 часа.
- Разработка процесса для ввода приложений в k8s, авторизации, рекомендуемых пайплайнов деплоя, etc.: N спринтов → поддержка.
- Апгрейд софта ноды (Docker, Kernel, etc.): несколько дней.
- Апгрейд версии кластера: не пробовали → квартал с двумя кластерами.

Предпосылки. Рост времени выполнения стандартных операций

Сейчас время выполнения Kubespray в sample-конфигурации:

```
root@b26f77349123:/kubespray# ansible-playbook -u ubuntu -b -i inventory/inventory.ini cluster.yml
```



Предпосылки. Чему мы научились и какие выводы сделали

Разобрались:

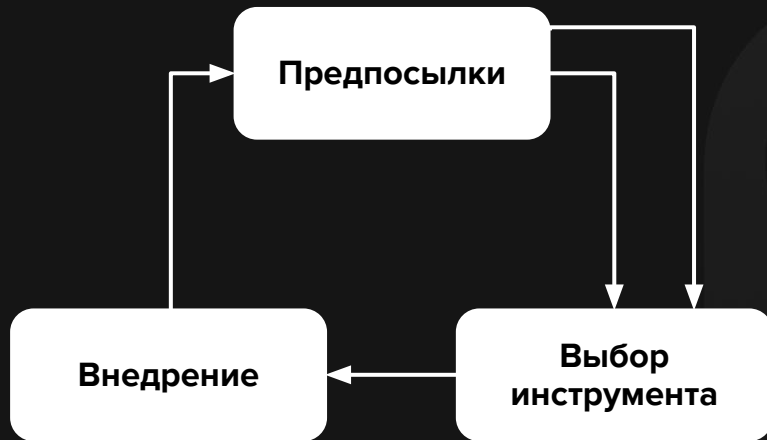
- В переменных Kubespray.
- Как работает CNI.
- Как работает k8s.

Какие выводы:

- Kubespray долго скипает задачи.
- Нужен не любой, а наш, настроенный, работающий кластер с нашими переменными.
- 3 года назад мастера через Kubespray скейлить нельзя, но мы примерно представляем как.
- Нужен инструмент, которому можно делегировать стандартные операции.

Выбор инструмента

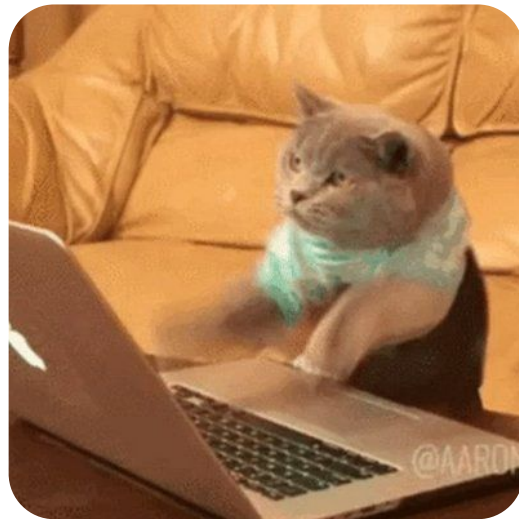
Q3 2019



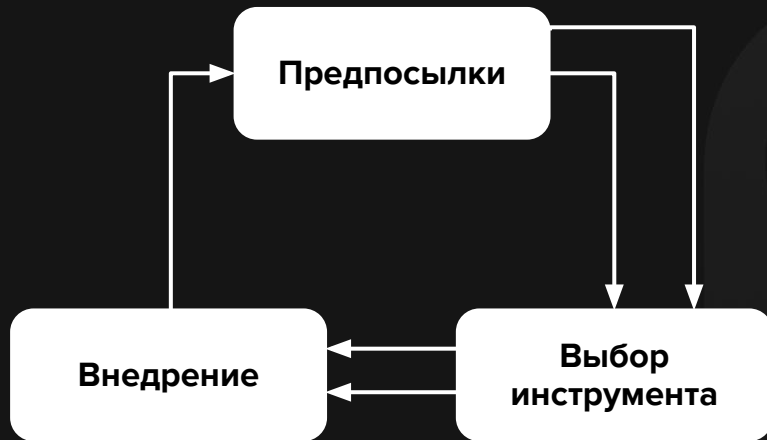
Выбор инструмента.

Kubeadm

- Вышел из беты.
- Команда получила контекст, поняли, что можем делегировать некоторые процессы.
- Kubeadm становится рекомендуемым способом оперирования k8s-кластером, в том числе и в официальной документации.



Внедрение. Q4 2019



Внедрение. Свой плейбук

```
# Create audit policy files on each master (for apiserver and falco)
- hosts: kube-master
  become: yes
  roles:
    - { role: kubernetes-audit-policy, tags: "kubernetes-audit-policy" }
  tags: ["setup-cluster", "k8s-audit-policy"]

# Setup first master (kubeadm init). It is executed only on first play on first master.
- hosts: kube-master[0]
  become: yes
  roles:
    - { role: kubeadm-init, tags: "kubeadm-init" }
  tags: ["setup-cluster", "kubeadm-init"]

# Generate join tokens and join new masters/nodes in cluster.
- hosts: kube-master[0]
  become: yes
  roles:
    - { role: kubeadm-join, tags: "kubeadm-join" }
  tags: ["setup-cluster", "kubeadm-join"]

# Setup calico using helm
- hosts: kube-master[0]
  become: yes
  roles:
    - { role: kubernetes-networking, tags: "kubernetes-networking" }
  tags: ["setup-cluster", "kubernetes-networking"]

# Add labels and annotations on nodes for manage taints using helm:
# app.kubernetes.io/managed-by: Helm
# meta.helm.sh/release-name: node-taints-labels
# meta.helm.sh/release-namespace: kube-system
- hosts: kube-master[0]
  become: yes
```

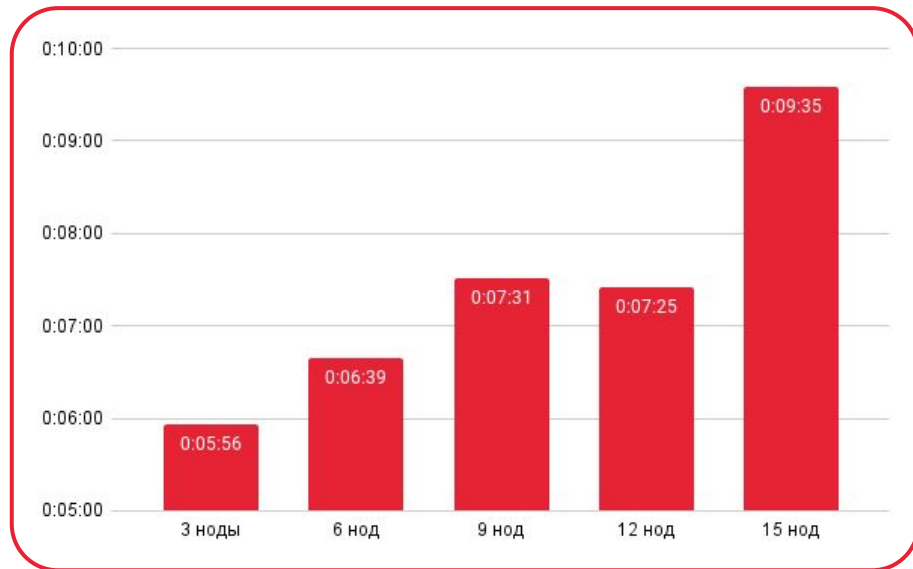
Что включила в себя эксплуатация

Стандартные операции:

- Раскатка кластера.
- Скейл кластера.
- Апгрейд софта ноды: Docker, Kernel, etc.
- Апгрейд версии кластера.
- Обновления сертификатов k8s, CNI, etc.

Размер кластера:

20 → 100 нод



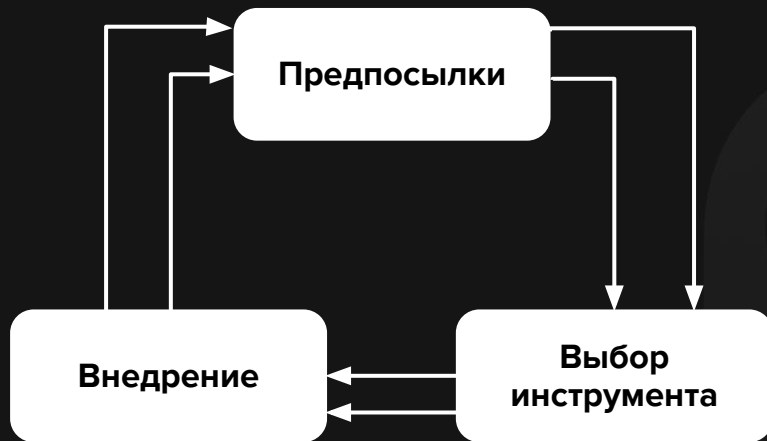
Что мы поняли еще через 2 года

Стандартные операции:

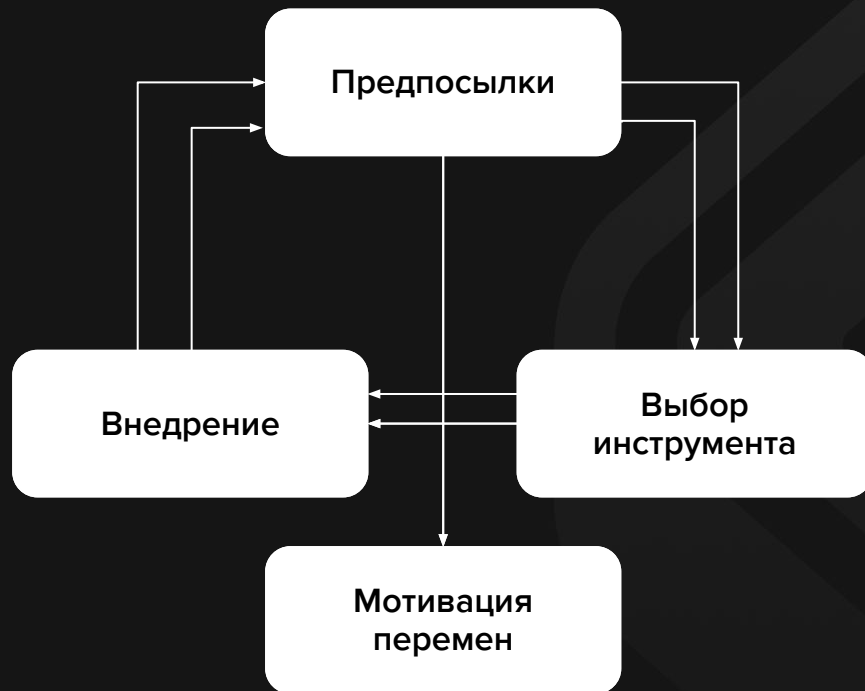
- Раскатка кластера: 10 минут → 1,5 часа.
- Скейл кластера: 10 минут → 1,5 часа.
- Апгрейд софта ноды (Docker, Kernel, etc) → несколько дней.
- Апгрейд версии кластера: пробовали → Kubeadm стал поддерживать апгрейд с 1.17.

Пришлось переезжать ещё раз в течение полугода.

Кажется, мы всё это проходили Q3 2021



Мотивация



Мотивация. Та же самая задача — развернуть кластер, но:

- Время поднятия < 3 минут на EC2.
- Число нод — произвольное.
- Число кластеров — произвольное.

Текущий флоу не подойдёт, потому что:

- Неудобно каждый раз прогонять плейбук на все тестовые среды, которых у нас ~ 200 в день.
- Мы не сможем сделать это руками даже с запуска джоб.
- Чтобы делать динамический инвентори, нужно делать автоматизацию.
- Если нужно делать автоматизацию, то зачем автоматизировать прокатку плейбука, если нам нужно автоматизировать поднятие кластера.

Мотивация. Profits

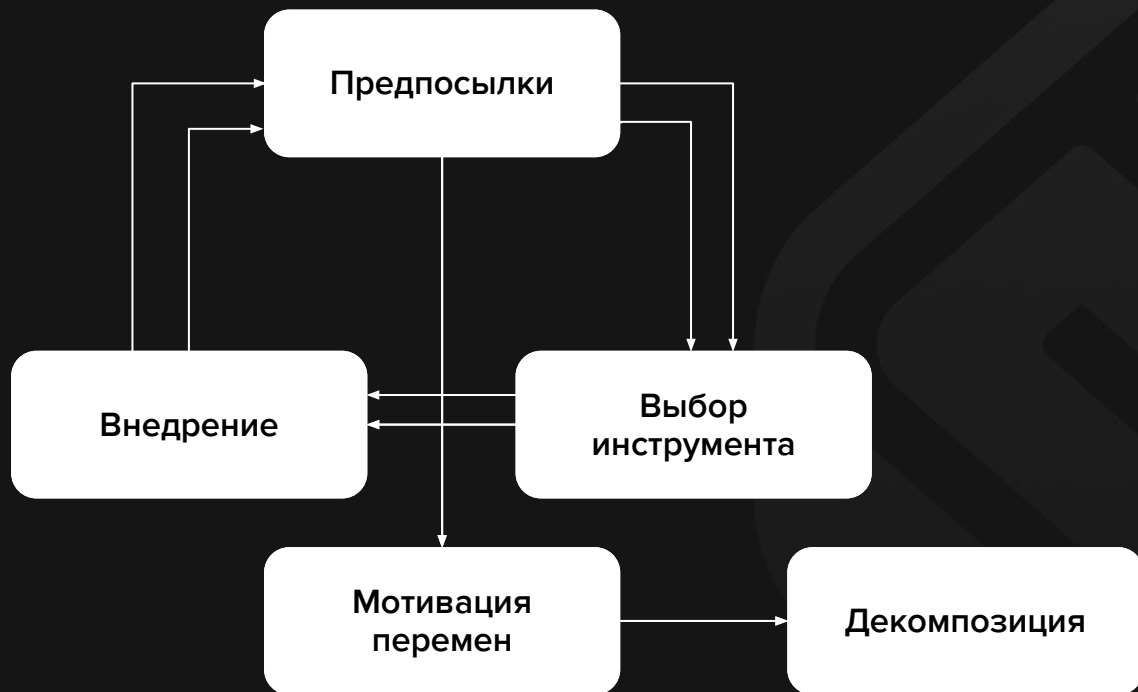
- Можно обеспечить постоянную и ежедневную проверку конфигурации на работающих тест-окружениях.
- Где-то здесь должно быть что-то, чтобы не катать наши железки 4 часа, попивая чай и ожидая stderr /dev/null от ансибла на одной из 100 нод.

```
fatal: [redacted] : FAILED! => {}
```

```
MSG:
```

```
Failed to get information on remote file (/root/.bashrc): Couldn't open /dev/null: Operation not permitted
```

Декомпозиция процесса



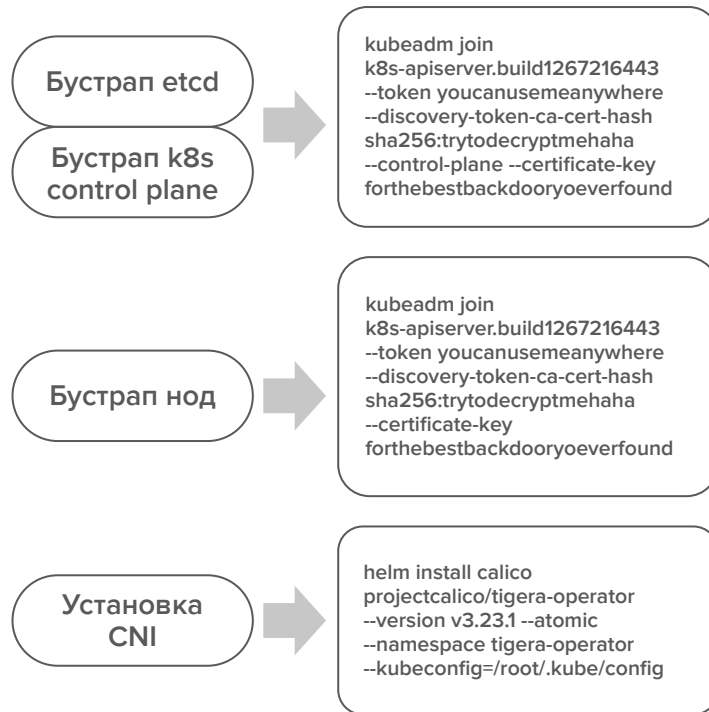
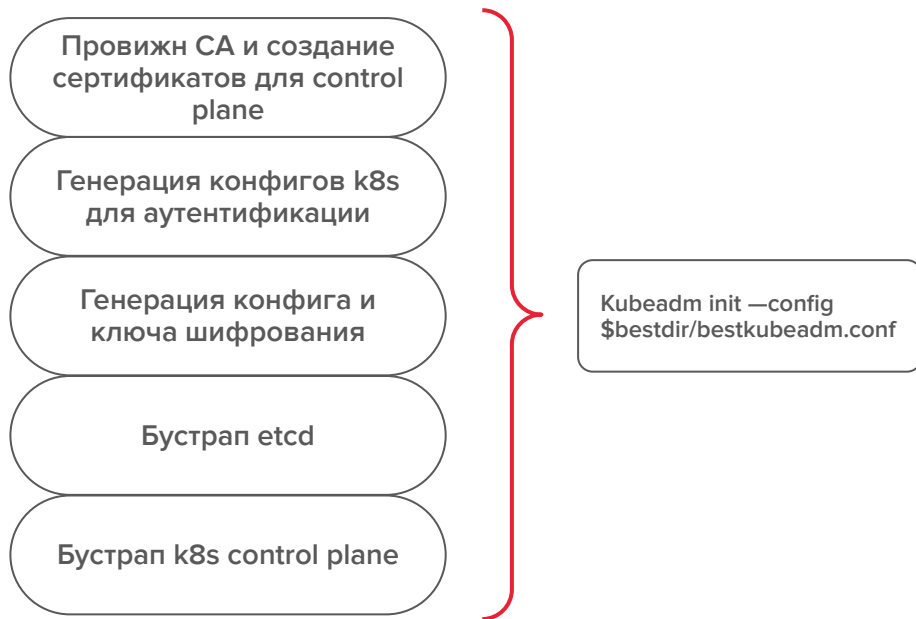
Декомпозиция процесса. Hard way

Вернёмся в самые дебри и посмотрим, что предлагает hard way:

Этапы по hard way	Число команд
Провижн CA и создание сертификатов для control plane	31
Генерация конфигов k8s для аутентификации	21
Генерация конфига и ключа шифрования	3
Бутстрап etcd	13 на каждом члене etcd-кластера
Бутстрап k8s control plane	33
Бутстрап нод	25
Конфигурация kubectl	5
Установка CNI	3
Установка DNS	6

Декомпозиция процесса. Документация k8s

Kubeadm



Декомпозиция процесса. Lifehack

А ноду забутстрапить можно и просто вот так:

Заранее сгенерируйте для kubelet'a сертификаты для авторизации

```
client-certificate: /var/lib/kubelet/pki/kubelet-client-current.pem  
client-key: /var/lib/kubelet/pki/kubelet-client-current.pem
```

Доставьте на ноду вместе с

```
"clientCAFile": "/etc/kubernetes/pki/ca.crt"  
/etc/kubernetes/kubelet/kubelet-config.json
```

И запустите kubelet с ключом --register-node

Прототипирование



Прототипирование. Делаем прототип

Ставим плейсхолдеры в groups vars, для kubeadm.conf

```
dex_endpoint: 'https://--DEX_DOMAIN--/'
kube_apiserver_fqdn: '--API_SERVER_DOMAIN--'
etcd_endpoints: ['--ETCD_DOMAIN--']

kubernetes_version: 1.21.13
kubernetes_utils_version: "{{ kubernetes_version }}-00"
```



```
root@localhost:/etc/kubernetes# cat kubeadm.conf
apiVersion: kubeadm.k8s.io/v1beta2
kind: ClusterConfiguration
imageRepository: k8s.gcr.io
certificatesDir: /etc/kubernetes/pki
clusterName: kubernetes
dns:
  type: CoreDNS
kubernetesVersion: "v1.21.13"
controlPlaneEndpoint: --API_SERVER_DOMAIN--:6443

etcd:
  local:
    serverCertSANS:
      - --ETCD_DOMAIN--
    peerCertSANS:
      - --ETCD_DOMAIN--
    dataDir: /var/lib/etcd
```

Прототипирование. Делаем прототип

Прописываем метаданные инстансам:

```
{
  "domainsEnv": "BASE_DOMAIN=*
API_SERVER_DOMAIN=* ETCD_DOMAIN=*",
  "masterAddr": "master01.build126.*",
  "kubenodesAddr": "node01.build126.*
node02.build126.* node03.build126.*",
  "myDomain": "node03.build126.*",
  "myRole": "node"
}
```

Current user data

[illegible]

Прототипирование. Делаем прототип

Пишем скрипт в cloud-init, чтобы забрать метаданные и запустить поднятие кластера:

```
function main () {
  log w "${FUNCNAME[0]}" "Starting Kubernetes cluster setup or join to existed"
  downloadMetadata
  extractMetadata
  exportGeneralUsageVars
  setHostname "${HOSTNAME}"
  case "${ROLE}" in
    "master" )
      setupMaster
      ;;

    "masterFollower" )
      setupMasterFollower
      ;;

    "node" )
      setupNode
      ;;

    * )
      log e "${FUNCNAME[0]}" "No instructions for role ${ROLE}"
      exit 1
  esac
}
```

Прототипирование. Что получилось

- Правда 3 минуты.
- Правда сколько угодно кластеров, пока есть инстансы в облаке.
- Правда непонятно, как скейлить и обновлять, но об этом позже. Мы же тут для тестовых коротко живущих сред поднимаем (ну, в общем случае).

Прототипирование. Как у вендоров?



EKS (join nodes)

```
root@ip-10-236-98-29 /]# grep eks /var/log/cloud-init-output.log
API_SERVER_URL=https://58F9FA156AC4C6840A53C91EA04A9AB9.gr7.eu-central-1.eks.amazonaws.com
/etc/eks/bootstrap.sh delete-me --kubelet-extra-args '--node-labels=eks.amazonaws.com/nodegroup-image=ami-004e93
```

```
TOKEN=$(get_token)
AWS_DEFAULT_REGION=$(get_meta_data 'latest/dynamic/instance-identity/document' | jq .region -r)
AWS_SERVICES_DOMAIN=$(get_meta_data 'latest/meta-data/services/domain')
MACHINE=$(uname -m)
PAUSE_CONTAINER_ACCOUNT=$(get_pause_container_account_for_region "${AWS_DEFAULT_REGION}")
PAUSE_CONTAINER_IMAGE=$(PAUSE_CONTAINER_IMAGE:$PAUSE_CONTAINER_ACCOUNT.dkr.ecr.${AWS_DEFAULT_REGION}.${AWS_SERVICES_DOMAIN}/eks/pause)
PAUSE_CONTAINER=${PAUSE_CONTAINER_IMAGE:$PAUSE_CONTAINER_VERSION}
CA_CERTIFICATE_DIRECTORY=/etc/kubernetes/pki
CA_CERTIFICATE_FILE_PATH=${CA_CERTIFICATE_DIRECTORY}/ca.crt
mkdir -p $CA_CERTIFICATE_DIRECTORY
aws eks wait cluster-active --region=${AWS_DEFAULT_REGION} --name=${CLUSTER_NAME}
aws eks describe-cluster --region=${AWS_DEFAULT_REGION} --name=${CLUSTER_NAME} --output=text --query 'cluster.certificateAuthorityData: certificateAuthority.data, endpoint: endpoint, serviceIpv4Cidr: kube-
netesNetworkConfig.serviceIpv4Cidr, serviceIpv6Cidr: kubernetesNetworkConfig.serviceIpv6Cidr, clusterIpFamily: kubernetesNetworkConfig.ipFamily, outpostArn: outpostConfig.outpostArns[0], id: id' > $DESCRIBE_CLUSTER_RESULT || rc=$?
B64_CLUSTER_CA=$(cat $DESCRIBE_CLUSTER_RESULT | awk '{print $1}')
K8SERVER_ENDPOINT=$(cat $DESCRIBE_CLUSTER_RESULT | awk '{print $3}')
CLUSTER_ID_IN_DESCRIBE_CLUSTER_RESULT=$(cat $DESCRIBE_CLUSTER_RESULT | awk '{print $4}')
OUTPOST_ARN=$(cat $DESCRIBE_CLUSTER_RESULT | awk '{print $5}')
SERVICE_IPV4_CIDR=$(cat $DESCRIBE_CLUSTER_RESULT | awk '{print $6}')
SERVICE_IPV6_CIDR=$(cat $DESCRIBE_CLUSTER_RESULT | awk '{print $7}')
echo $B64_CLUSTER_CA | base64 -d > $CA_CERTIFICATE_FILE_PATH
sed -i s,MASTER_ENDPOINT,$K8SERVER_ENDPOINT,g /var/lib/kubelet/kubeconfig
sed -i s,AWS_REGION,$AWS_DEFAULT_REGION,g /var/lib/kubelet/kubeconfig
DOMAIN_NAME=$(echo "$K8SERVER_ENDPOINT" | awk -F/ '{print $3}') | awk -F: '{print $1}'
getent hosts "$DOMAIN_NAME" | shuf >> /etc/hosts
sed -i s,CLUSTER_NAME,$CLUSTER_ID,g /var/lib/kubelet/kubeconfig
KUBELET_EXTRA_ARGS="--bootstrap-kubeconfig /var/lib/kubelet/kubeconfig SKUBELET_EXTRA_ARGS"
## Kubelet service configuration
MCO=$(get_meta_data 'latest/meta-data/network/interfaces/macs/' | head -n 1 | sed 's/./\//')
DNS_CLUSTER_IP=$(DNS_CLUSTER_IP)
KUBELET_CONFIG=/etc/kubernetes/kubelet/kubelet-config.json
echo "${jq 'clusterDNS+=("${DNS_CLUSTER_IP}")' $KUBELET_CONFIG}" > $KUBELET_CONFIG
MAX_PODS_FILE=/etc/eks/eni-max-pods.txt
set +o pipefail
MAX_PODS=$(cat $MAX_PODS_FILE | awk '"/$(INSTANCE_TYPE)-unset/"/ { print $2 }')
MAX_PODS=$(/etc/eks/max-pods-calculator.sh --instance-type-from-ids --cni-version 1.10.0 --show-max-allowed)
echo "${jq '++ {\"evictionHard\": {\"memory.available\": \"100Mi\", \"nodefs.available\": \"10%\", \"nodefs.inodesFree\": \"5%\"}}' $KUBELET_CONFIG}" > $KUBELET_CONFIG
echo "${jq --arg mibibytes_to_reserve \"$mibibytes_to_reserve\"M --arg cpu_millicores_to_reserve \"$cpu_millicores_to_reserve\"m '++ {\"kubeReserved\": {\"cpu\": $cpu_millicores_to_reserve, \"ephemeral-storage\": \"1Gi\", \"memory\": $mibibytes_to_reserve}}' $KUBELET_CONFIG}" > $KUBELET_CONFIG
echo "${jq '++ {\"maxPods\": $MAX_PODS}' $KUBELET_CONFIG}" > $KUBELET_CONFIG
mkdir -p /etc/systemd/system/kubelet.service.d
cat <<EOF > /etc/systemd/system/kubelet.service.d/10-kubelet-args.conf
cat <<EOF > /etc/systemd/system/containerd.service.d/10-compat-symlink.conf
mkdir -p /etc/docker
bash -e "/sbin/iptables-save > /etc/sysconfig/iptables"
cp -v /etc/eks/iptables-restore.service /etc/systemd/system/iptables-restore.service
sudo chown root:root /etc/systemd/system/iptables-restore.service
systemctl daemon-reload
systemctl enable iptables-restore
echo "$DOCKER_CONFIG_JSON" > /etc/docker/daemon.json
echo "${jq 'bridge= \"docker0\"' | 'live-restore=false' $docker/docker.json}" > /etc/docker/daemon.json
systemctl enable kubelet
systemctl start kubelet
```

Прототипирование. Как у вендоров?

Digital ocean

```
root@pool-1l5m5oomo-7pgeo:/var/lib/cloud/scripts/per-instance# find /var/lib/cloud/ -name *k8s*  
/var/lib/cloud/scripts/per-instance/000-k8saas
```

```
K8SAAS_PUBLIC_ADDRESS="$(curl --retry 30 -fSs http://169.254.169.254/metadata/v1/interfaces/public/0/ipv4/address)"  
K8SAAS_PRIVATE_ADDRESS="$(curl --retry 30 -fSs http://169.254.169.254/metadata/v1/interfaces/private/0/ipv4/address)"  
export K8SAAS_PUBLIC_ADDRESS  
export K8SAAS_PRIVATE_ADDRESS  
K8SAAS_PROVIDER_ID="digitalocean://$(curl --retry 30 -sSf http://169.254.169.254/metadata/v1/id)"  
export K8SAAS_PROVIDER_ID  
done < <(grep ^k8saas_ < /var/lib/cloud/instance/user-data.txt)  
  
case "${K8SAAS_ROLE:-}" in  
    master)  
        rm -f /etc/systemd/system/kubelet.service.d/10-kubeadm.conf  
        rm -f /lib/systemd/system/kubelet.service  
        /opt/k8saas/bootstrap-master  
    ;;  
    kubelet)  
        rm -f /etc/systemd/system/kubelet.service.d/10-kubeadm.conf  
        rm -f /lib/systemd/system/kubelet.service  
        /opt/k8saas/bootstrap-kubelet  
    ;;  
    *)  
        echo 'Assuming kubeadm usage'  
        exit 0  
    ;;  
    **)  
        echo "Don't know role ${K8SAAS_ROLE}. Valid options: [master/kubelet]"  
        exit 1  
    ;;  
esac  
~
```

Прототипирование. Как у вендоров?

Digital ocean

[illegible]

Прототипирование. Как у вендоров?

А что это за сервисы такие красивые в GKE

```
gke-my-first-cluster-1-default-pool-3a2323fe-f5m5 ~ # grep -r kube /var/log/cloud-init-output.log
Created symlink /etc/systemd/system/kubernetes.target.wants/kube-node-installation.service → /etc/systemd/system/kube-node-installation.service.
Created symlink /etc/systemd/system/kubernetes.target.wants/kube-node-configuration.service → /etc/systemd/system/kube-node-configuration.service.
Created symlink /etc/systemd/system/kubernetes.target.wants/kube-container-runtime-monitor.service → /etc/systemd/system/kube-container-runtime-monitor.service.
Created symlink /etc/systemd/system/kubernetes.target.wants/kubelet-monitor.service → /etc/systemd/system/kubelet-monitor.service.
Created symlink /etc/systemd/system/kubernetes.target.wants/kube-logrotate.timer → /etc/systemd/system/kube-logrotate.timer.
Created symlink /etc/systemd/system/kubernetes.target.wants/kube-logrotate.service → /etc/systemd/system/kube-logrotate.service.
Created symlink /etc/systemd/system/multi-user.target.wants/kubernetes.target → /etc/systemd/system/kubernetes.target.
```



```
gke-my-first-cluster-1-default-pool-3a2323fe-f5m5 ~ # systemctl list-dependencies | grep kube
● kubernetes.target
● └─ kube-container-runtime-monitor.service
○ └─ kube-logrotate.service
○ └─ kube-logrotate.timer
● └─ kube-node-configuration.service
● └─ kube-node-installation.service
● └─ kubelet-monitor.service
```

Прототипирование. Как у вендоров?

GKE — главное

```
gke-my-first-cluster-1-default-pool-3a2323fe-f5m5 ~ # cat /home/kubernetes/bin/configure-helper.sh | wc -l
3725
gke-my-first-cluster-1-default-pool-3a2323fe-f5m5 ~ # cat /home/kubernetes/bin/configure.sh | wc -l
1269
```



Метаданные всё так же забираются со своего API

```
function download-kube-master-certs {
# Fetch kube-master-certs from GCE metadata server.
(
umask 077
local -r tmp_kube_master_certs="/tmp/kube-master-certs.yaml"
# shellcheck disable=SC2086
retry-forever 10 curl ${CURL_FLAGS} \
-H "X-Google-Metadata-Request: True" \
-o "${tmp_kube_master_certs}" \
http://metadata.google.internal/computeMetadata/v1/instance/attributes/kube-master-certs
# Convert the yaml format file into a shell-style file.
eval "$(python3 -c ''
```

Прототипирование. Как у вендоров?



Мы уже крутые? Не совсем

Что имеем:

- HW-кластера с их проблемами всё так же на месте.
- У нас появилась ещё одна точка поддержки.

Какие сделали выводы:

- По-другому можно.
- По-другому несложно.
- Возможно проще, чем было раньше.

Мы уже крутые? Не совсем



Внедрение. Чего минимально не хватает каждой системе

Hardware:

- Система инвентаризации.
- Сервер метаданных.
- Пулл-система обработки метаданных.

Внедрить можно за месяц,
согласовывать нужно 2 года.

Для облака:

- Процесс скейлинга кластера.
- Процесс апгрейда существующих кластеров.

Внедрить несложно, сложно
дождаться, когда станет нужно
для бизнеса.





Выводы. Какие проблемы можно приобрести

- Это уже не пуш-система. Сделанная вручную ошибка скорее всего уже будет подхвачена — это типичная проблема пулл-систем.
- Fail fast — must have.
- Нужно решить вопрос A/B-тестирования конфигурации.
- Система становится сложнее:

Было	Стало
Добавил в инвентори — пошёл смотреть в лог.	Несколько независимых систем: <ul style="list-style-type: none">• provisioner инстансов;• поставщик метаданных;• поставщик снапшотов;• система, оперирующая нодой на основе метаданных.

Выводы

Над нами сейчас вот такие задачи, и именно таким подходом мы хотим их решить:

- Апгрейднуть ядер для 100 нод за час.
- Уверенность в disaster recovery.
- Все окружения раскатываются одинаково.
- Нет зависимости от размера или количества кластеров.

Выводы

В каком порядке мы будем идти дальше:

- Добавить скейл и апгрейд на EC2-кластера.
- Понять все доп. системы, необходимые вокруг Bare Metal.
- Согласовать с другими командами и бизнесом инструментарий.
- Внедрить.
- Перезагрузить ноды кластера.

Выводы. Оглядываясь назад на 5,5 лет



Thank you!

Обратная связь
и комментарии
по докладу по ссылке

Илья Устинов
📍 @Crowd292



HighLoad ++
2022

